# Mono-Forward: Backpropagation-Free Algorithm for Efficient Neural Network Training Harnessing Local Errors

**James Gong, Bruce Li**
The University of Auckland
Auckland, New Zealand
{hgon777,tli389}@aucklanduni.ac.nz

**Waleed Abdulla**
The University of Auckland
Auckland, New Zealand
w.abdulla@auckland.ac.nz

## Abstract

Backpropagation is the standard method for achieving state-of-the-art accuracy in neural networks training, but it often imposes high memory costs and lacks biological plausibility. In this paper, we introduce the Mono-Forward algorithm, a purely local layerwise learning method inspired by Hinton's Forward-Forward framework. Unlike backpropagation, Mono-Forward optimizes each layer solely with locally available information, eliminating the reliance on global error signals. We evaluated Mono-Forward on multi-layer perceptrons and convolutional neural networks across multiple benchmark, including MNIST, Fashion-MNIST, CIFAR-10, and CIFAR-100. The test results show that Mono-Forward consistently matches or surpasses the accuracy of backpropagation across all tasks, with significantly reduced and more even memory usage, better parallelizability, and comparable convergence rate.

## 1 Introduction

Backpropagation [1], while foundational in training neural networks [2], faces critical limitations in both deep learning and neuroscience, highlighting the importance of exploring alternative methodologies.

The concept of Backward Locking exemplifies a significant bottleneck inherent to BP, where weight updates across the network must await the completion of both forward and backward passes for each data batch, hampering the efficient distribution of computation and parallelization across the network [3–5].

In BP, error gradients can exhibit significant variations in magnitude as they are propagated backwards through the network's layers, leading to two prominent issues: vanishing and exploding gradients. Vanishing gradients occur when the gradients diminish to such small values that they fail to effectively update the weights of earlier layers, thus severely hindering the training of deep neural networks. On the other hand, exploding gradients present a challenge by causing disproportionately large updates to the weights, potentially destabilizing the network. The implementation of solutions such as those documented in [6–8] to mitigate these issues inevitably introduces additional complexity into both the network architecture and the training algorithm.

The biological implausibility of BP is highlighted by its dependence on a global error signal to update the entire network, a mechanism that lacks substantial support in biological neural systems. According to [9], there is minimal evidence that such a global error feedback mechanism exists within the neural architecture of the brain. Neuroscience research generally supports the notion that the cortex does not engage in the backward propagation of errors across its layers, which is a fundamental aspect of the BP algorithm. This view was notably advanced in [10], which argued against the presence of such error propagation mechanisms in biological networks. This discrepancy highlights a critical divergence between artificial neural networks and biological systems, raising questions about the applicability of BP-based methods to neurological models. Instead, local learning paradigms are more closely aligned with local plasticity rules like Spike-Timing-Dependent Plasticity (STDP) [11–13] in biological neural circuits, potentially making them more relevant for understanding real brain function and informing the design of more biologically inspired artificial learning algorithms.

This divergence has spurred interest in alternative learning algorithms that may be more congruent with biological mechanisms, such as feedback alignment [14], direct feedback alignment [15], and the Forward-Forward algorithm [16]. Despite the advancements in alternative methodologies, BP, taking the advantages of global error, remains the algorithm for delivering the highest accuracy [14, 16–24], especially in complex datasets such as CIFAR-100 [25].

In this study, we introduce the Mono-Forward algorithm, a layer-wise greedy backpropagation-free training method inspired by the Forward-Forward algorithm, which demonstates accuracy comparable to or exceeding that of backpropagation on the benchmark datasets with significantly reduced memory consumption, better parallelizability, equivalent prediction efficiency, and comparable convergence rate.

## 2    The Forward-Forward Algorithm

In December 2022, Hinton proposed the Forward-Forward (FF) algorithm as an alternative to BP. FF replaces the forward and backward passes in BP by two forward passes, termed the positive and negative forward passes. This approach aims to optimize neural networks in a layer-wise, greedy manner without global backward error propagation. Each layer adjusts its parameters based on the "goodness" scores, evaluated separately in the positive and negative passes.

In FF, the positive goodness, evaluated during the positive pass, is computed as the sum of squares of the neuron activations for each layer $i$, with correctly labeled one-hot encoded data as the input. Similarly, the negative goodness, obtained during the negative pass, is computed with incorrectly labeled one-hot encoding data as the input. The goal is to adjust the network parameters such that positive goodness exceeds a specific threshold while pushing negative goodness well below the threshold. Formally, the goodness values of a particular layer $i$ is

$$G_{pos,i} \triangleq \sum_{n'=0}^{n'=n-1} (A_{i,n'}^2), \quad x \to x_{pos} \tag{1}$$

$$G_{neg,i} \triangleq \sum_{n'=0}^{n'=n-1} (A_{i,n'}^2), \quad x \to x_{neg} \tag{2}$$

where $G_{pos,i}$ and $G_{neg,i}$ represent the layer-wise positive and negative goodness respectively. $A_{i,n'}$ represents the activation of the $n'_{th}$ neuron, where there are $n$ neurons in total. $x_{pos}$ and $x_{neg}$ indicate positive and negative data input.

The positive and negative loss values for each layer $i$ are computed as follows, with threshold parameter $\theta$:

$$L_{pos,i} \triangleq log(1 + e^{\theta - G_{pos,i}}), \quad L_{neg,i} \triangleq log(1 + e^{G_{neg,i} - \theta}) \tag{3}$$

$$L_{total,i} \triangleq \frac{1}{2}(L_{pos,i} + L_{neg,i}) \tag{4}$$

The learning process of the network involves calculating the partial derivatives of the total loss function with respect to each weight at the current layer $i$, and this error information is detached from the latter layers. Then each weight parameter is adjusted using gradient descent.

In the prediction phase of FF, the input data are encoded with each possible label, generating distinct datasets that are fed into the network sequentially. For each dataset, the goodness values from all layers are summed to compute a total goodness score, and the label with the highest score is selected as the output.

The Forward-Forward algorithm is often regarded as more biologically plausible than backpropagation because it avoids the propagation of the global error signal. However, the need for a "negative pass" challenges this plausibility [26–28]. Also, in Section 5 of the FF paper [16], our experimental findings align with those reported by Hinton, demonstrating that effective learning with FF necessitates pairing the positive and negative passes within the same epoch. This observation challenges the hypothesis that the negative pass could be analogous to the biological brain's sleep phase, thus casting doubt on the biological plausibility of the negative pass in the FF framework.

Another efficiency issue with FF is its requirement for $m$ data passes to make a prediction, where $m$ represents the number of categories. This multi-pass approach becomes particularly inefficient with a large number of categories. To mitigate this inefficiency, Hinton introduced a one-pass version of the FF algorithm. However, this adaptation results in a trade-off, reducing accuracy. Additionally, the FF algorithm faces challenges in achieving high convergence rate due to the absence of a direct connection between labels and input data.

## 3 The Mono-Forward Algorithm

Mono-Forward (MF) is a greedy, layer-wise training method that optimizes neural networks through a single forward pass, significantly enhancing its computational efficiency compared to FF. Unlike FF, MF requires only one pass for predictions, regardless of the number of categories. It establishes a direct connection between labels and input data, which results in a more stable and faster convergence.

### 3.1 Training in MF

MF employs a specialized approach to calculate the "goodness" score at each layer, which plays a central role in optimizing the network. Each layer in MF contains a projection matrix $M_i$ whose dimensions are $m \times n$, where $m$ represents the number of possible categories, and $n$ is the number of neurons in layer $i$. The goodness score $G_i$ for the $i_{\text{th}}$ layer is computed as follows:

$$G_i \triangleq a_i \times M_i^\top \tag{5}$$

In this equation, $a_i$ denotes the activation values of the neurons at layer $i$, and $M_i$ is the matrix specific to that layer. The transposed matrix $M_i^T$ allows the multiplication to align correctly, yielding a matrix where each row in $G_i$ corresponds to the goodness scores across all categories for each example in the batch. Thus, each element of $G_i$ represents a category-specific goodness score for each instance, facilitating a direct comparison across categories.

For training, MF leverages these goodness scores in conjunction with the cross-entropy loss function, which is defined as follows:

$$\mathcal{L}_i \triangleq - \sum_{c=1}^{m} y_c \log(\sigma(G_{ic})) \tag{6}$$

Here, $y_c$ is a binary indicator (0 or 1) if a class label $c$ is the correct classification for observation $i$, and $\sigma(G_{ic})$ is the softmax function applied to the goodness scores $G_i$ for each class, normalizing them into a probability distribution over the classes. The cross-entropy loss thus aims to maximize the predicted probability of the correct label while minimizing that of the incorrect labels. The loss value, $\mathcal{L}_i$, evaluated is then used to update the layer weights and the projection matrix specific to that layer. The update of the parameters is dictated by the equations (7) and (8), where $\eta$ represents the learning rate of the algorithm, $\mathbf{a_i}$ is the activation vector in the layer $i$, and $\mathbf{z_i}$ is the linear output before activation.

$$\mathbf{W}_i \leftarrow \mathbf{W}_i - \eta\left(\frac{\partial \mathcal{L}_i}{\partial \mathbf{G}_i} \cdot \frac{\partial \mathbf{G_i}}{\partial \mathbf{a_i}} \cdot \frac{\partial \mathbf{a_i}}{\partial \mathbf{z_i}} \cdot \frac{\partial \mathbf{z_i}}{\partial \mathbf{W_i}}\right) \tag{7}$$

$$\mathbf{M}_i \leftarrow \mathbf{M}_i - \eta\left(\frac{\partial \mathcal{L}_i}{\partial \mathbf{G}_i} \cdot \frac{\partial \mathbf{G_i}}{\partial \mathbf{M_i}}\right) \tag{8}$$

In this configuration, each layer produces activations, $a_i$, that can be effectively used by the projection matrix to maximize the distinction between the goodness scores of the correct label and those of the incorrect labels. These activations from the current layer $i$ serve as input to the subsequent layer. Consequently, each successive layer is equipped to learn from the features that have been optimized in the previous layer for maximum discriminatory power between the true and false labels in a greedy and feed-forward manner.

The following algorithm summarizes the training process for a single batch using the Mono-Forward method, incorporating the computation of goodness scores and the application of the cross-entropy loss for network optimization:

---
**Algorithm 1** Layer-Wise Training for One Batch

---
1: **Input:** $\mathbf{X}_{\text{batch}}$, $\mathbf{y}_{\text{batch}}$        ▷ Batch data, labels
2: $\mathbf{a}_0 \leftarrow \mathbf{X}_{\text{batch}}$        ▷ Initialize activations with input batch
3: **for** $i \leftarrow 1,$ num_layers **do**
4:      $\mathbf{z}_i \leftarrow \mathbf{a}_{i-1}\mathbf{W}_i$        ▷ Optional Bias Term
5:      $\mathbf{a}_i \leftarrow \phi(\mathbf{z}_i)$        ▷ Activation
6:      $\mathbf{G}_i \leftarrow \mathbf{a}_i\mathbf{M}_i^\top$
7:      $\mathcal{L}_i \leftarrow \text{CrossEntropy}(\mathbf{y}_{\text{batch}}, \mathbf{G}_i)$
8:      $\mathbf{W}_i \leftarrow \mathbf{W}_i - \eta\frac{\partial \mathcal{L}_i}{\partial \mathbf{W}_i}$
9:      $\mathbf{M}_i \leftarrow \mathbf{M}_i - \eta\frac{\partial \mathcal{L}_i}{\partial \mathbf{M}_i}$
10:     Pass $\mathbf{a}_i$ to the next layer
11: **end for**

---

Overall, this approach enables the network to learn from a single forward pass by converting the implicit label-input connections found in FF into explicit connections. In FF, labels are embedded as part of the input and there is no explicit guidance in the algorithm on the specialness of the label, but in MF, the goodness score for each label is calculated using a distinct set of parameters in $M_i$, enabling explicit label-input connections in the algorithm. Moreover, this method allows the algorithm to learn distinctions between the correct label and all incorrect labels simultaneously. In contrast, the FF algorithm compares the correct label with only one incorrect label at a time, limiting its scope of comparative learning.

Another perspective on understanding the MF is through the concept of "future projection". In this view, each projection matrix in the layers acts as a condensed projection of the future operations, akin to the form of mini neural networks. This projection guides the current layer to anticipate the operations that will be applied to the activations, and then it allows the current layer to "prepare" its activations to better align with these anticipated future operations, effectively tailoring its output to fit this "future projection" while ensuring consistent learning across layers. In this algorithm, the projection matrix is treated as the standard form of "future projection".

### 3.2 Prediction in MF

In MF, predictions can be made using two distinct approaches: one inherited from FF methodology, and the other resembling traditional BP techniques. The effectiveness of each method may vary slightly depending on the specific task at hand, with one approach potentially yielding slightly higher accuracy than the other for different types of data or learning objectives.

In FF approach, each layer of the network computes a set of goodness scores corresponding to each potential label. For example, at layer 0, the tensor $G_0$ contains goodness scores, with each of the $m$ entries representing a different label, where $m$ is the total number of possible labels. These goodness scores from each layer are subsequently aggregated to determine the final goodness scores across the network. The label that corresponds to the highest aggregated goodness score is then selected as the final prediction.

In BP approach, the goodness scores are computed solely by the last layer of the network. This setup is similar to traditional BP methods, where the output layer directly follows the last hidden layer. Each score represents a potential label, and the label corresponding to the highest goodness score is immediately selected as the prediction. This direct method simplifies the prediction process by consolidating the decision-making to the final layer.

In comparison to FF approach, the BP approach in the MF framework requires fewer network parameters. This is because the projection matrices in the preceding layers, which are essential in the FF method, are unnecessary in the BP approach and can be omitted during the prediction phase. Consequently, this streamlined BP method achieves the identical level of prediction efficiency as traditional BP, being faster than the one-pass version of FF, which was introduced by Hinton as a less optimal version designed to enhance prediction speed but at the expense of reduced accuracy.

### 3.3 Parallelizability, Transparency, and Hot-Plugging

One of the limitations of BP is its inherent sequential nature, requiring the completion of the entire forward pass before beginning the backward pass. The network must store activations and intermediate values during the forward pass to compute gradients in the backward pass. This sequential dependency inhibits full parallelization, potentially reducing computational efficiency and lengthening training times for large models. In contrast, MF allows for weight adjustments as soon as the input has passed through a single layer, significantly reducing the time delay compared to BP, where weight adjustments must wait until the entire network has completed the forward pass.

Another feature of BP is that this algorithm prioritizes the adjustment of weights in the final layer of a neural network, with modifications to earlier layers made solely to support these end-layer changes, as dictated by the chain rule. This method, while effective for enhancing training accuracy, does not necessarily promote uniform learning across all layers, as exemplified by the gradient vanishing issue. This lack of consistent learning can potentially lead to overfitting as the final layers in BP can compensate for earlier layers' shortcomings, especially in scenarios where there are many more parameters than necessary for achieving full training accuracy–a common situation in both artificial and biological neural networks. However in MF, every layer must contribute meaningfully to class discrimination, preventing the network from "putting all its complexity" into the final layers and thus encouraging better generalization.

Furthermore, MF provides transparency in understanding the contribution of each layer to the overall network performance. Unlike traditional approaches, each layer in MF is capable of making independent predictions, as demonstrated in Table 1. This characteristic allows us to evaluate the performance of each layer individually.

4

In the BP prediction mode of MF, only the last layer contributes to the final prediction. By analyzing outputs from the network with three hidden layers in Table. 1, it becomes evident that the third hidden layer's performance is lower than that of the second layer. This suggests that the third layer is over-extracting features, leading to overfitting. In this example, MF makes it possible to identify such issues without requiring additional training of a two-layer network for comparison, as is typically necessary in BP to achieve optimal performance.

Table 1: Comparison of classification accuracy (mean $\pm$ standard deviation) on CIFAR-10 for Backpropagation (BP) and the Mono-Forward (MF) algorithm using locally connected multilayer perceptron (MLP) architectures. The MLP processes 3 input channels and uses kernel sizes of 4, 8, and 16 for configurations with 3 hidden layers, while kernel sizes of 4 and 8 are used for configurations with 2 hidden layers. Each result is averaged over at least 3 different random seeds.

| Algorithm | Total Hidden Layers | Accuracy (%) |
|---|---|---|
| BP | 2 Hidden Layers | $59.40 \pm 0.13$ |
| MF (FF Pred) | 2 Hidden Layers | $61.61 \pm 0.06$ |
| MF Layer 1 | 2 Hidden Layers | $59.57 \pm 0.21$ |
| MF Layer 2 | 2 Hidden Layers | $\mathbf{61.71 \pm 0.18}$ |
| BP | 3 Hidden Layers | $58.79 \pm 0.07$ |
| MF (FF Pred) | 3 Hidden Layers | $\mathbf{61.78 \pm 0.03}$ |
| MF Layer 1 | 3 Hidden Layers | $59.57 \pm 0.21$ |
| MF Layer 2 | 3 Hidden Layers | $61.71 \pm 0.18$ |
| MF Layer 3 | 3 Hidden Layers | $60.11 \pm 0.36$ |

In addition, MF naturally allows seamless "plugging" or "unplugging" of layers without requiring any modifications to preceding layers. The outputs of the first two layers in a trained two-layer network are identical to those of the first two layers in a trained three-layer network. This consistency enables the earlier layers to act as independent modules, effectively "black-boxing" their computations, offering simplified weight management and modular design. However, BP is not inherently modular and requires manual intervention, such as freezing weights, to emulate the modularity and independence seen in MF.

Another notable advantage of this modularity is its ability to reduce the overall loss when part of the network is damaged or compromised. If a specific layer is affected by damage—whether due to model corruption, hardware failure, or adversarial perturbation—only the damaged layer needs to be retrained. The unaffected preceding layers retain their functionality and do not require recalibration, preserving the network's overall structure. Furthermore, even without any retraining, the network can still make predictions based on the outputs of earlier layers, albeit with reduced accuracy.

This behavior mirrors observations in stroke patients, who, despite experiencing damage to specific brain regions, are often able to perform tasks with diminished accuracy or efficiency. For instance, patients with damage to motor control areas of the brain may still retain the ability to move but often exhibit slower or less precise movements [29]. Similarly, individuals with damage to language processing regions, such as the Broca area, may still communicate but with reduced fluency or articulation [30].

In biological neural networks, dendrites perform localized non-linear computations, acting as "mini-neural networks" that process and integrate synaptic inputs before transmitting signals to the soma [31–34]. Similarly, in the Mono-Forward (MF) algorithm, each layer's projection matrix functions as an independent computational unit, transforming activations into goodness scores for class predictions. These matrices refine inputs locally, akin to how dendritic subunits process inputs independently, contributing to the overall function of the neuron. By treating each layer as a self-contained "mini-network," MF mirrors the decentralized, hierarchical nature of dendritic computation, further enhancing the biological plausibility of the algorithm.

## 4   Experiments

To evaluate the feature learning capabilities of MF, both Multi-Layer Perceptrons (MLPs) and Convolutional Neural Networks (CNNs) were employed. The network architectures were designed to match the complexity of each dataset. For the MNIST [35] and Fashion-MNIST [36] datasets, a 2x1000 ReLU MLP was trained using all algorithms listed in Table 2. In contrast, for the more complex CIFAR-10 [37] and CIFAR-100 [25] datasets, a 3x2000 ReLU MLP architecture was employed.

For CNN experiments, a convolutional network with four convolutional layers was used. Each convolutional layer applies a kernel size of $3 \times 3$ with padding to preserve spatial dimensions, followed by ReLU activation and average pooling to reduce feature map size. Specifically, the architecture consists of 64, 128, 256, and 512 feature maps in the respective layers, culminating in a fully connected layer of size $512 \times 2 \times 2$, which maps to the output classes. All networks are optimized with Adaptive Moment Estimation (Adam) [38]. Hyperparameters, such as learning rate and batch size, were fine-tuned based on testing set performance.

In MLP experiments, the performance of MF was evaluated using the FF prediction and BP prediction approach, along with other algorithms such as FF, BP, Feedback Alignment (FA), and Direct Feedback Alignment (DFA). The results, summarized in Table 3, reveal that while the other non-backpropagation algorithms fail to achieve the same level of accuracy as BP when applied to increasingly complex datasets, MF consistently outperforms BP in terms of accuracy, demonstrating its effectiveness across all evaluated datasets.

| Algorithm | MNIST (%) | F-MNIST (%) | CIFAR-10 (%) | CIFAR-100 (%) |
|---|---|---|---|---|
| BP | $99.52 \pm 0.05$ | $92.63 \pm 0.06$ | $77.80 \pm 0.49$ | $42.10 \pm 0.03$ |
| MF (FF Pred) | $\mathbf{99.58} \pm 0.01$ | $\mathbf{93.98} \pm 0.08$ | $\mathbf{82.39} \pm 0.25$ | $\mathbf{54.77} \pm 0.25$ |
| MF (BP Pred) | $99.58 \pm 0.04$ | $93.47 \pm 0.16$ | $81.89 \pm 0.28$ | $51.26 \pm 0.28$ |
| DFA | $98.53 \pm 0.17$ | $87.87 \pm 0.21$ | $57.53 \pm 0.19$ | $24.54 \pm 0.60$ |
| FA | $98.99 \pm 0.11$ | $90.04 \pm 0.93$ | $62.90 \pm 0.61$ | $26.02 \pm 0.82$ |

Table 2: Comparative accuracy (mean $\pm$ standard deviation) of BP, MF, DFA, and FA algorithms on different datasets. Results are obtained using a vanilla CNN with network configuration of 64, 128, 256, and 512 neurons per layer. Each experiment is conducted with at least 3 different random seeds. The best result in each dataset is highlighted in bold.

In the CNN experiments, the performance of MF was compared with the same set of algorithms, excluding FF. This is because FF is not intended for the convolutional neural networks as mentioned by Hinton in his original paper [16]. The results, presented in Table 2, show that MF consistently outperforms BP in all data sets. In particular, the performance gap between MF and BP increases as the complexity of the data set increases. Moreover, MF demonstrates a greater advantage over BP in CNNs compared to MLP, whereas the other non-backpropagation algorithms seem to adapt less well in CNNs.

| Algorithm | MNIST (%) | F-MNIST (%) | CIFAR-10 (%) | CIFAR-100 (%) |
|---|---|---|---|---|
| BP | $98.69 \pm 0.05$ | $90.27 \pm 0.19$ | $54.25 \pm 0.26$ | $27.64 \pm 0.28$ |
| FA | $92.71 \pm 0.05$ | $84.49 \pm 0.15$ | $51.70 \pm 0.35$ | $17.01 \pm 0.20$ |
| DFA | $98.32 \pm 0.08$ | $89.17 \pm 0.07$ | $53.86 \pm 0.35$ | $21.76 \pm 0.28$ |
| FF | $97.35 \pm 0.07$ | $87.88 \pm 0.13$ | $46.55 \pm 0.22$ | - |
| MF (FF Pred) | $\mathbf{98.74} \pm 0.01$ | $\mathbf{90.52} \pm 0.02$ | $\mathbf{56.99} \pm 0.17$ | $\mathbf{29.05} \pm 0.13$ |
| MF (BP Pred) | $\mathbf{98.74} \pm 0.02$ | $90.51 \pm 0.05$ | $56.40 \pm 0.20$ | $28.27 \pm 0.25$ |

Table 3: Comparative accuracy (mean $\pm$ standard deviation) of BP, FA, DFA, FF, and MF algorithms on different datasets. Results are obtained using an MLP architecture with $2 \times 1000$ layers for MNIST and F-MNIST, and $3 \times 2000$ layers for CIFAR-10 and CIFAR-100. Each experiment is conducted with at least 3 different random seeds. The best result achieved in each dataset is highlighted in bold. FF fails to converge on CIFAR-100.

In the optimal form of the Forward-Forward algorithm, it takes $m$ data passes to make a single prediction, where $m$ is the number of categories. To address this limitation, Hinton proposed the one-pass FF that requires only a single data pass to make predictions in all occasions, however, this version is described as the suboptimal version of FF and comes with a cost in accuracy. Although some forms of FF, such as Symba FF [39], and [24], are proposed to improve accuracy or convergence, many of them still fall short in prediction efficiency compared to BP, which is intrinsic in the design of FF. In our algorithm, despite being inspired by FF, we do not utilize the concept of "positive" or "negative" data, and therefore demands only one data pass in both FF prediction and BP prediction approach in the algorithm's most natural form. To assess the prediction efficiency of MF, both of the prediction approaches in MF are set up to compare against FF, one-pass FF, and BP on MNIST dataset. The results, shown in Table 4, highlight that while MF with FF prediction is slower than BP, MF with BP prediction achieves the same prediction time complexity and parameter count as BP. This level of efficiency is unattainable even in the one-pass version of FF, demonstrating MF's ability to balance performance and efficiency effectively. In particular, both prediction mechanisms in MF consistently outperform BP in the experiments we have conducted in Tables 3 and 2.

| Algorithm | Single Time (ms) | Batch Time (ms) | Accuracy (%) | Parameter Number |
|---|---|---|---|---|
| FF | 5.6 | 13.2 | 97.33 | **1,784k** |
| One-Pass FF | 0.5 | 1.8 | 95.54 | 1,804k |
| BP | **0.3** | **1.1** | 98.75 | 1,794k |
| MF (FF Pred) | 1.7 | 3.1 | 98.75 | 1,804k |
| MF (BP Pred) | **0.3** | **1.1** | **98.76** | 1,794k |

Table 4: Comparison of prediction times and accuracy for FF, BP, and MF algorithms on a network architecture with $2 \times 1000$ ReLU neurons, evaluated on an NVIDIA 4090 GPU and MNIST dataset. Results are measured in milliseconds and averaged over 100 runs. "Batch" indicates the prediction time required for a batch of 10,000 images, while "Single" represents the prediction time for a single image.

The training process for MF involves updating both layer weights and projection weights, along with a forward pass of the data. In contrast, BP only updates the layer weights and performs a forward pass. Although MF requires more operations than BP in total, its better parallelizability results in a shorter wall time per epoch. This is reflected in the experiment results shown in Table 5, where the training time per epoch is consistently lower for MF compared to BP, both for full batch and mini-batch training on MNIST and CIFAR-10 datasets [37].

| Dataset | Batch Type | Network | Mean Time (s) | Std (s) | Speed Ratio |
|---|---|---|---|---|---|
| MNIST | Mini Batch | MF | 0.71 | 0.10 | 2.28 |
| | | BP | 1.62 | 0.21 | – |
| | Full Batch | MF | 0.14 | 0.03 | 1.57 |
| | | BP | 0.22 | 0.07 | – |
| CIFAR10 | Mini Batch | MF | 0.66 | 0.11 | 3.33 |
| | | BP | 2.20 | 0.26 | – |
| | Full Batch | MF | 0.04 | 0.02 | 1.57 |
| | | BP | 0.06 | 0.01 | – |

Table 5: Training Time Comparison between MF and BP. Mean epoch time and standard deviation for a 2-hidden-layer network, each layer consisting of 1000 ReLU neurons, with either full batch size or a mini-batch size of 64. Results are averaged over 10 epochs, with experiments conducted on a Nvidia 4080 Laptop GPU. The speed ratio is calculated as the ratio of BP time to MF time (BP time / MF time).

One advantage of local learning algorithms over backpropagation, which relies on a global error signal, is reduced memory usage [40, 41]. In layerwise learning algorithms, memory requirements depend primarily on the activations and computations of the current layer. To verify that MF retains this benefit, we analyzed the memory usage during training for both MF and BP, as shown in Fig. 1 and 2. The results confirm that BP requires more memory than MF. In addition, MF exhibits a narrower gap between memory peaks and valleys compared to BP, indicating a more consistent memory consumption and better utilization of hardware resources. Moreover, even memory usage can minimize delays caused by memory reallocation or swapping, contributing to smoother and faster training processes [42, 43].

To illustrate the relationship between memory usage and the number of layers, we plotted the peak memory allocated during training against the number of layers in the network, as shown in Fig. 3. The results compare BP and MF across both MLP and CNN architectures. Specifically for the MLP architecture, where each additional layer contains a fixed size of 1000 neurons, memory usage increases linearly with the number of layers. In Fig. 3a, the slopes of the best-fit lines were computed: BP has a slope of 204.00, while MF has a significantly lower slope of 16.16. Notably, similar to other layerwise training algorithms, if layer weights were persisted to storage instead of being retained in memory, the memory usage would theoretically remain constant, independent of the depth of the network.

To evaluate the convergence speed of MF, we tested a 4-layer MLP network, each with 1000 neurons, using the Adam optimizer. Both algorithms were evaluated under identical conditions, with a batch size of 256 and a learning rate of 0.001. The experiment was run at least 10 times with different seeds and the average test error was recorded. The resulting convergence curves are shown in Fig. 4. Due to the layerwise nature of MF, the training dynamics of the network remain unaffected by its depth. This characteristic is evident in Fig. 5, which highlights the contrasting
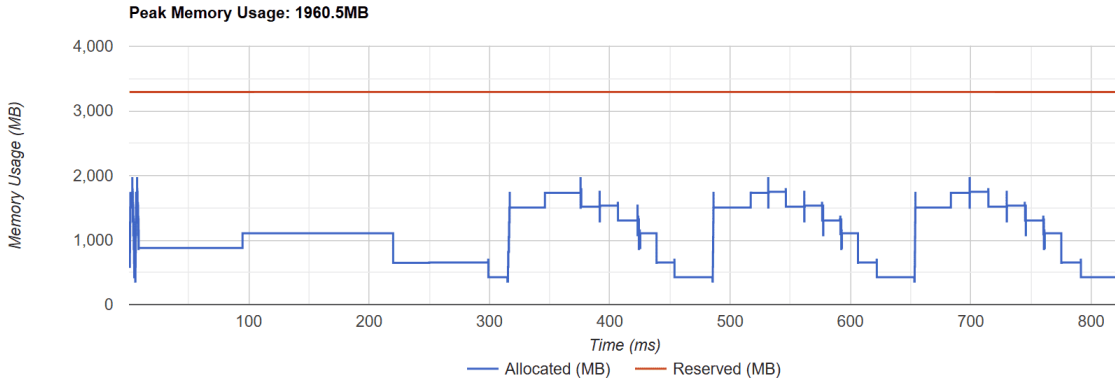
Figure 1: Memory Consumed during Training under BP. This experiment utilizes MNIST dataset on a network of size 5 * 2000 with batch size 30000.
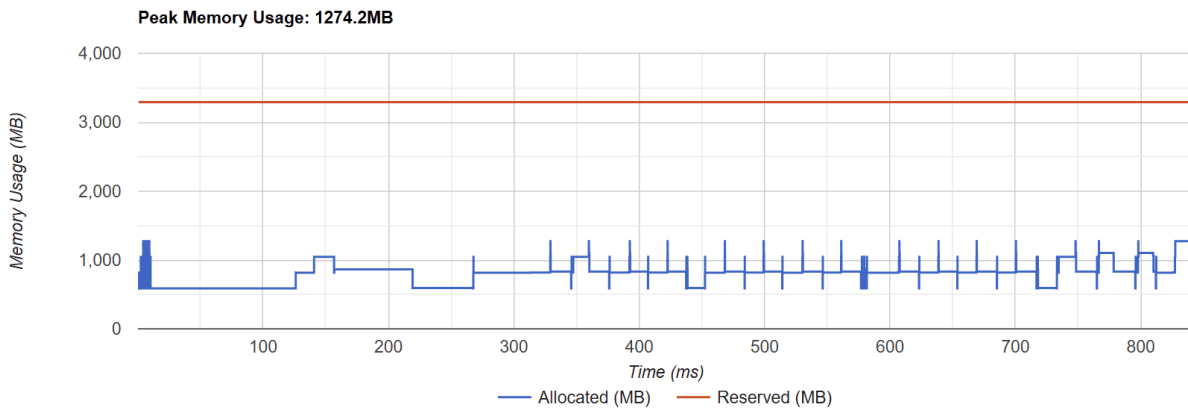


Figure 2: Memory Consumed during Training under MF. This experiment utilizes MNIST dataset on a network of size 5 * 2000 with batch size 30000.
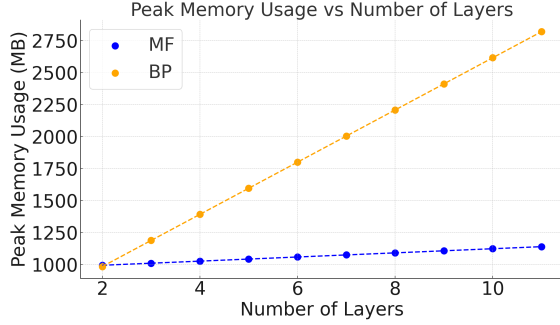
behavior of BP and MF as the number of layers increases. Specifically, BP demonstrates a significant degradation in convergence rate with deeper architectures (Fig. 5a), whereas MF maintains consistent convergence rate regardless of depth (Fig. 5b).

Lastly, to evaluate the effectiveness of MF beyond image classification tasks, we compared MF with BP on tabular data and text sentiment analysis using the Transaction Prediction dataset [44], Breast Cancer dataset [45], and IMDB dataset [46]. The results are summarized in Table 6.
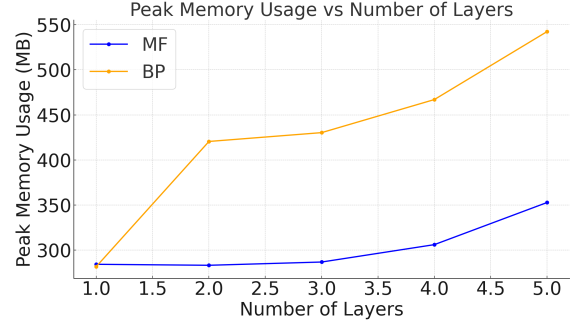
| Dataset | Model | Results (%) |
|---|---|---|
| Transaction Prediction | BP | $91.17 \pm 0.03$ |
| | MF | $91.17 \pm 0.02$ |
| Breast Cancer | BP | $97.37 \pm 0.00$ |
| | MF | $97.37 \pm 0.00$ |
| IMDB | BP | $87.41 \pm 0.17$ |
| | MF | $88.37 \pm 0.15$ |

Table 6: Performance comparison between BP and MF. Transaction Prediction is evaluated using an MLP with three hidden layers of 200 neurons each, Breast Cancer using an MLP with two hidden layers of 100 neurons each, and IMDB using an MLP with an embedding layer of 50 neurons and a single hidden layer of 100 neurons. These architectures were selected for optimal performance on their respective tasks. Each configuraion is at least tested with 3 different seeds, and the means and standard deviations are recorded.
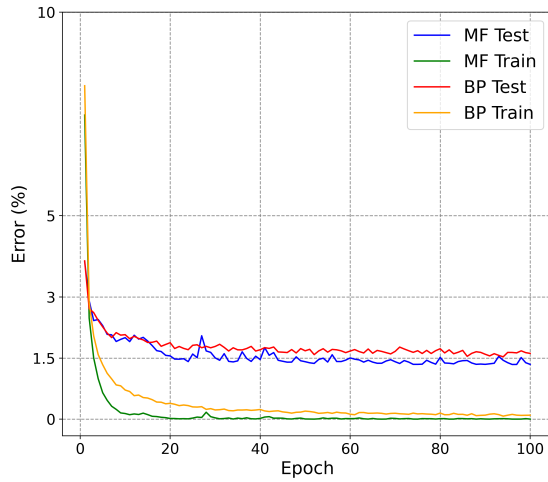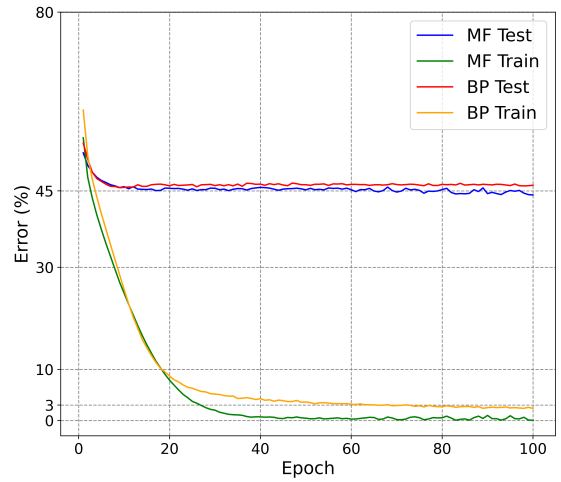
8

(a) MLP Architecture



(b) CNN Architecture

Figure 3: Memory Comparison between Backpropagation (BP) and Mono-Forward (MF) during Training. For the MLP architecture, the experiment involves varying numbers of 1000 ReLU neuron layers on the MNIST dataset using a full batch size. For the CNN architecture, the setup includes different number of convolutional layers with 64, 128, 256, 512, and 512 neurons on the CIFAR-10 dataset with a batch size of 256. The peak memory usage during training is recorded.
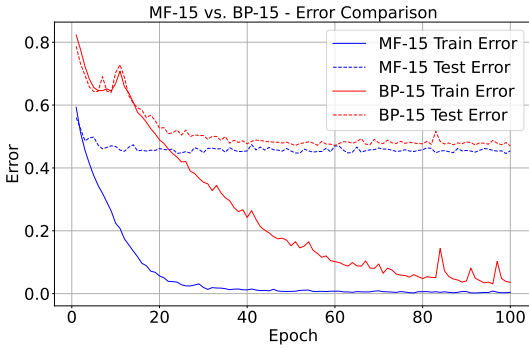


(a) MNIST training and testing error curve



(b) CIFAR training and testing error curve

Figure 4: Comparison of Convergence Rates between MF and BP. The experiments were conducted using a MLP network with $4 \times 1000$ neurons, optimized using the Adam optimizer with a learning rate of 0.001. Results are averaged over at least 10 different random seeds for each dataset.
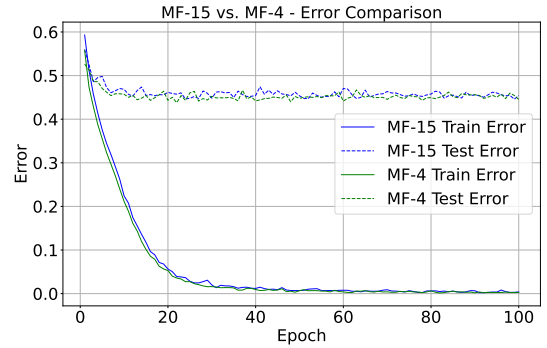
## 5  Conclusion

The experimental results demonstrate that, compared to BP, both MF prediction modes achieve higher accuracy in the testing set while offering significantly reduced and more consistent memory consumption and shorter training times per epoch. This is achieved with a convergence rate and prediction efficiency (BP prediction mode only) comparable to those of BP, using only a single forward pass and local errors. Furthermore, MF's mechanism promotes transparency and facilitates dynamic modifications to neural network architectures. It also enhances parallelizability and ensures that training process is not inhibited by network depth, making it a robust and scalable alternative to traditional backpropagation.

In contrast to FF, MF eliminates the need for one-hot encoding and negative data, which are less biologically plausible [26–28]. MF computes error derivatives for all false classes in a single forward pass using cross-entropy loss, avoiding FF's reliance on squared outputs and threshold tuning. Furthermore, MF requires only one data pass for training and prediction, offering significantly improved prediction efficiency over FF. In addition, the projection matrix mechanism, which treats each layer as a mini-neural network, closely mimics the biological concept of dendritic computation, where individual dendrites process localized inputs to contribute to the overall neural response.

(a) 15 Layers BP vs 15 Layers MF

(b) 15 Layers MF vs 4 Layers MF

Figure 5: Comparison of Convergence Rates between MF and BP. The experiments were conducted using a MLP network with $n \times 1000$ neurons, optimized using the Adam optimizer with a learning rate of 0.001, with batch size 256. MF-15 means a 15 layers MLP network trained with MF.

Last but not the least, while it is widely acknowledged that local learning algorithms, which rely solely on local information, offer advantages in memory efficiency and biological plausibility, they typically achieve lower testing accuracy compared to backpropagation, which leverages global information. MF introduces a novel and biologically plausible approach that not only retains the advantages of local learning algorithms but also surpasses backpropagation in accuracy.

# References

[1] David E. Rumelhart, Geoffrey E. Hinton, and Ronald J. Williams. Learning representations by back-propagating errors. *Nature*, 323(6088):533–536, October 1986.

[2] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–444, May 2015.

[3] Zhouyuan Huo, Bin Gu, qian Yang, and Heng Huang. Decoupled parallel backpropagation with convergence guarantee. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 2098–2106. PMLR, 10–15 Jul 2018.

[4] Wojciech Czarnecki, Grzegorz Swirszcz, Max Jaderberg, Simon Osindero, Oriol Vinyals, and Koray Kavukcuoglu. Understanding synthetic gradients and decoupled neural interfaces. 03 2017.

[5] Max Jaderberg, Wojciech Marian Czarnecki, Simon Osindero, Oriol Vinyals, Alex Graves, David Silver, and Koray Kavukcuoglu. Decoupled neural interfaces using synthetic gradients. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 1627–1635. PMLR, 06–11 Aug 2017.

[6] Xavier Glorot, Antoine Bordes, and Yoshua Bengio. Deep sparse rectifier neural networks. In Geoffrey Gordon, David Dunson, and Miroslav Dudík, editors, *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, volume 15 of *Proceedings of Machine Learning Research*, pages 315–323, Fort Lauderdale, FL, USA, 11–13 Apr 2011. PMLR.

[7] Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural networks. In Sanjoy Dasgupta and David McAllester, editors, *Proceedings of the 30th International Conference on Machine Learning*, volume 28 of *Proceedings of Machine Learning Research*, pages 1310–1318, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR.

[8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. pages 770–778, 06 2016.

[9] Yuhang Song, Thomas Lukasiewicz, Zhenghua Xu, and Rafal Bogacz. Can the brain do backpropagation? -exact implementation of backpropagation in predictive coding networks. *Advances in Neural Information Processing Systems*, 33:22566–22579, January 2020.

[10] Francis Crick. The recent excitement about neural networks. *Nature*, 337(6203):129–132, January 1989.

[11] Guoqiang Bi and Mu-ming Poo. Synaptic modifications in cultured hippocampal neurons: Dependence on spike timing, synaptic strength, and postsynaptic cell type. *The Journal of neuroscience : the official journal of the Society for Neuroscience*, 18:10464–72, 01 1999.

[12] Sen Song, Kenneth Miller, and L.F. Abbott. Competitive hebbian learning through spike timing-dependent plasticity. *Nature neuroscience*, 3:919–26, 10 2000.

[13] Natalia Caporale and Yang Dan. Spike timing–dependent plasticity: A hebbian learning rule. *Annual review of neuroscience*, 31:25–46, 02 2008.

[14] Maria Refinetti, Stéphane d'Ascoli, Ruben Ohana, and Sebastian Goldt. Align, then memorise: the dynamics of learning with feedback alignment, 2021.

[15] Arild Nøkland. Direct feedback alignment provides learning in deep neural networks, 2016.

[16] Geoffrey Hinton. The forward-forward algorithm: Some preliminary investigations, 2022.

[17] Brian Crafton, Abhinav Parihar, Evan Gebhardt, and Arijit Raychowdhury. Direct feedback alignment with sparse connections for local learning, 2019.

[18] Julien Launay, Iacopo Poli, and Florent Krzakala. Principled training of neural networks with direct feedback alignment. *ArXiv*, abs/1906.04554, 2019.

[19] Benjamin Scellier and Y. Bengio. Equilibrium propagation: Bridging the gap between energy-based models and backpropagation. *Frontiers in Computational Neuroscience*, 11, 05 2017.

[20] Maxence Ernoult, Julie Grollier, Damien Querlioz, Yoshua Bengio, and Benjamin Scellier. Updates of equilibrium prop match gradients of backprop through time in an rnn with static input, 2019.

[21] Dong-Hyun Lee, Saizheng Zhang, Asja Fischer, and Yoshua Bengio. Difference target propagation, 2015.

[22] Erkki Oja. Simplified neuron model as a principal component analyzer. *Journal of Mathematical Biology*, 15(3):267–273, November 1982.

[23] Peter Földiák. Forming sparse representations by local anti-hebbian learning. *Biol Cybern*, 64:165–70, 02 1990.

[24] Guy Lorberbom, Itai Gat, Yossi Adi, Alexander Schwing, and Tamir Hazan. Layer collaboration in the forward-forward algorithm. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(13):14141–14148, Mar. 2024.

[25] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-100 (canadian institute for advanced research).

[26] Demis Hassabis, Dharshan Kumaran, Christopher Summerfield, and Matthew Botvinick. Neuroscience-inspired artificial intelligence. *Neuron*, 95:245–258, 07 2017.

[27] Blake A Richards, Timothy P Lillicrap, Philippe Beaudoin, Yoshua Bengio, Rafal Bogacz, Amelia Christensen, Claudia Clopath, Rui Ponte Costa, Archy de Berker, Surya Ganguli, Colleen J Gillon, Danijar Hafner, Adam Kepecs, Nikolaus Kriegeskorte, Peter Latham, Grace W Lindsay, Kenneth D Miller, Richard Naud, Christopher C Pack, Panayiota Poirazi, Pieter Roelfsema, João Sacramento, Andrew Saxe, Benjamin Scellier, Anna C Schapiro, Walter Senn, Greg Wayne, Daniel Yamins, Friedemann Zenke, Joel Zylberberg, Denis Therien, and Konrad P Kording. A deep learning framework for neuroscience. *Nature Neuroscience*, 22(11):1761–1770, October 2019.

[28] Randall C. O'Reilly. Six principles for biologically based computational models of cortical cognition. *Trends in Cognitive Sciences*, 2(11):455–462, 1998.

[29] John Krakauer, S. Thomas Carmichael, Dale Corbett, and George Wittenberg. Getting neurorehabilitation right: What can be learned from animal models? *Neurorehabilitation and neural repair*, 26:923–31, 03 2012.

[30] Martha Taylor Sarno. Assessment of aphasia and related disorders. *Physical Therapy*, 53(2):225–226, 02 1973.

[31] Alon Poleg-Polsky, Bartlett Mel, and Jackie Schiller. Polsky, a., mel, b.w. & schiller, j. computational subunits in thin dendrites of pyramidal cells. nat. neurosci. 7, 621-627. *Nature neuroscience*, 7:621–7, 07 2004.

[32] Nelson Spruston. Spruston n. pyramidal neurons: dendritic structure and synaptic integration. nat rev neurosci 9: 206-221. *Nature reviews. Neuroscience*, 9:206–21, 04 2008.

[33] Christof Koch. *Biophysics of Computation: Information Processing in Single Neurons*. Oxford University Press, 11 1998.

[34] Michael London and Michael Häusser. Dendritic computation. *Annual review of neuroscience*, 28:503–32, 2005.

[35] Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE Signal Processing Magazine*, 29(6):141–142, 2012.

[36] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms, 2017. cite arxiv:1708.07747Comment: Dataset is freely available at https://github.com/zalandoresearch/fashion-mnist Benchmark is available at http://fashion-mnist.s3-website.eu-central-1.amazonaws.com/.

[37] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-10 (canadian institute for advanced research).

[38] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

[39] Heung-Chang Lee and Jeonggeun Song. Symba: Symmetric backpropagation-free contrastive learning with forward-forward algorithm for optimizing convergence, 2023.

[40] Marc'Aurelio Ranzato, Christopher Poultney, Sumit Chopra, and Yann Lecun. Efficient learning of sparse representations with an energy-based model. 01 2006.

[41] Geoffrey E. Hinton, Simon Osindero, and Yee-Whye Teh. A fast learning algorithm for deep belief nets. *Neural Computation*, 18(7):1527–1554, 2006.

[42] Minsoo Rhu, Natalia Gimelshein, Jason Clemons, Arslan Zulfiqar, and Stephen W. Keckler. vdnn: virtualized deep neural networks for scalable, memory-efficient neural network design. In *The 49th Annual IEEE/ACM International Symposium on Microarchitecture*, MICRO-49. IEEE Press, 2016.

[43] NVIDIA. Gpu memory management, 2020. NVIDIA Developer Documentation.

[44] Santander. Santander customer transaction prediction, 2018.

[45] D. Dua and C. Graff. Uci machine learning repository: Breast cancer wisconsin (diagnostic) data set, 2017.

[46] Andrew L Maas, Raymond E Daly, Peter T Pham, Dan Huang, Andrew Y Ng, and Christopher Potts. Learning word vectors for sentiment analysis. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 142–150, 2011.